

# Percolation and epidemics in clustered networks

Joel C. Miller\*

*Harvard School of Public Health, Boston, MA  
Fogarty Institute, National Institutes of Health, Bethesda, MD*

(Dated: June 15, 2009)

The social networks that infectious diseases spread along are typically clustered. Because of the close relation between percolation and epidemic spread, the behavior of percolation in such networks gives insight into infectious disease dynamics. A number of authors have studied clustered networks, but the networks often contain preferential mixing between high degree nodes. We introduce a class of random clustered networks and another class of random unclustered networks with the same preferential mixing. Percolation in the clustered networks reduces the component sizes and increases the epidemic threshold compared to the unclustered networks.

Classical random networks contain few short cycles, and the proportion of nodes in short cycles goes to zero as the number of nodes increases. In contrast social networks typically contain many short cycles. We refer to such networks as *clustered* networks. The impact of clustering on percolation properties is usually difficult to calculate because cycles prevent the use of branching process arguments, but it is widely expected that clustering significantly alters percolation.

Typically studies of infectious disease spread assume that outbreaks begin with a single infected node. The disease travels to each susceptible neighbor independently with probability  $T$ , the *transmissibility*, and the node recovers. The process repeats. We focus on diseases for which recovery provides immunity, so recovered nodes are not susceptible. Typically the outbreak dies out stochastically or becomes an epidemic and spreads until the number of susceptible nodes is reduced.

It is well-established that for fixed  $T$ , epidemic spread can be mapped to a bond percolation problem wherein each edge is kept with probability  $T$  [5, 8, 10, 12, 13, 16]. If we perform percolation on the network and then choose the initial infection, the disease spreads from that initial infection along edges of the percolated network, and so an epidemic occurs iff the initial node is in the giant component. The size of the epidemic matches the size of the giant component. This establishes that the probability and fraction infected in epidemics are equal if  $T$  is fixed and all edges are independent [22].

Because social networks frequently exhibit clustering, a number of studies have investigated the impact of clustering on epidemic problems [1, 4, 7, 9, 14, 17, 20, 21]. Some have found that clustering reduces the sizes of epidemics and raises the epidemic threshold. That is, clustering reduces the size of giant components and raises the percolation threshold. However, others have shown that clustering appears to reduce the threshold. Consequently epidemics would be possible at lower transmissibility in the presence of clustering.

This discrepancy occurs because there are many ways

used to generate clustered networks. It is difficult to separate the impact of clustering from other features introduced by the network generation process.

In this article we introduce a new algorithm to generate random clustered networks [23]. The clustered networks have correlations between degrees in a well-defined manner which can lead to assortativity, the tendency for nodes to contact nodes of similar degree. We show how to generate unclustered networks with the same correlations. We can make analytic comparisons between the two, and so clearly separate the effect of clustering from degree correlations. We show that although the clustered networks can have a reduced threshold compared to purely random networks of the same degree distribution, that is entirely an artifact of the assortativity. Compared to an unclustered network of the same degree correlations, the clustered networks result in smaller epidemics and higher epidemic threshold.

This article is organized as follows: we first introduce our clustered and unclustered networks. We then calculate and compare the epidemiological quantity  $\mathcal{R}_0$  which measures how many new infections a typical infected node causes. Finally, we calculate the final size/probability of epidemics assuming constant  $T$ .

## I. THE NETWORKS

We model our approach after standard algorithms for Configuration Model (CM) networks [3, 15, 18]. CM networks are useful because all edges from a node are independent of one another, in the sense that whether an epidemic results from following one edge is independent of the result along any other edge because short cycles are negligible.

### A. Clustered Networks

We begin with  $N$  nodes. To each node  $u$  we assign two degrees, an *independent edge* degree  $k_I$  and a triangle degree  $k_\Delta$ . The joint probability of  $k_I$  and  $k_\Delta$  is given by  $p(k_I, k_\Delta)$ . Then  $u$  will be part of  $k_\Delta$  triangles and have  $k_I$  other edges. Each triangle and edge from  $u$  will

---

\*Electronic address: joel.c.miller.research@gmail.com

be independent of other triangles and edges in the same way that edges in CM networks are independent.

We create an independent stub list and a triangle stub list. We place  $u$  into the independent stub list  $k_I$  times and into the triangle stub list  $k_\Delta$  times. Once all nodes are placed into the lists, we randomize them. We then take the pairs of nodes in positions  $2n$  and  $2n + 1$  of the independent list and join them, and the triples in positions  $3n$ ,  $3n + 1$ , and  $3n + 2$  and join them into a triangle. Some repeated edges or loops or short cycles other than the triangles we impose may appear, but their impact is negligible as  $N \rightarrow \infty$  [24].

This algorithm inevitably segregates those nodes with a high proportion of triangles from those nodes with a low proportion of triangles. If the degrees of nodes with many triangles differ from the degrees of nodes with few triangles, then this effect will cause correlation of different degrees. In order to isolate the impact of clustering, we must be able to compare percolation in these clustered networks with percolation in networks whose nodes are segregated in the same way.

## B. Unclustered, Segregated Networks

For comparative purposes we develop a corresponding unclustered network with the same segregation as the clustered networks. Given the joint distribution  $p(k_I, k_\Delta)$  of independent and triangle degrees, we create a new network where nodes are assigned *blue* and *red* degrees such that  $k_b = k_I$  and  $k_r = 2k_\Delta$ . The joint distribution is given by  $p_u(k_b, k_r) = p(k_b, k_r/2)$ .

We proceed as before. We create a blue and red list, and pair nodes in positions  $2n$  and  $2n + 1$  in the blue list and then repeat with the red list, joining pairs, not triples. The resulting network has the same segregation as the corresponding clustered network, but short cycles are negligible.

## II. $\mathcal{R}_0$

$\mathcal{R}_0$  is usually defined as the number of new infections caused by an average infected individual. Occasionally alternate definitions are used, but in some way it represents the number of new infections attributed to an average infected individual.  $\mathcal{R}_0 = 1$  is the threshold below which epidemics have zero probability (*i.e.*, the percolated network has no giant component). If  $\mathcal{R}_0 > 1$  then epidemics are possible, but not guaranteed.

### A. Clustered Networks

To simplify the analysis, first assume that  $u$ ,  $v$ , and  $w$  are members of a triangle and  $u$  becomes infectious first. There are multiple ways that both  $v$  and  $w$  can become infected from edges within the triangle, but they all have

the same impact on the epidemic. It is convenient to treat infections of  $v$  and  $w$  as if they came from from  $u$  regardless of the actual path followed.

Thus if  $u$  becomes infected, then with probability  $2T^2(1-T) + T^2 = 3T^2 - 2T^3$  it is credited with infecting both  $v$  and  $w$ , and with probability  $2T(1-T)^2$  it is credited with infecting just 1. With probability  $(1-T)^2$  it infects neither. Thus the expected number of infections per triangle is  $2T(1+T-T^2)$ . In spirit this approach is similar to that of [2]. For book-keeping purposes, we define the rank  $s$  of a node as follows: the index case has rank 0. Each node  $v$  is then assigned rank  $s$  to be the shortest path of infectious contacts from the index case to  $v$ , under the rule above for crediting infections.

This allows us to define a  $2 \times 2$  next-generation matrix [6]. We separate those nodes infected along an independent edge from those nodes infected along a triangle edge [25]. We define  $c_{II}$  and  $c_{\Delta I}$  to be the number of infections that a node infected from an independent edge is expected to cause along independent and triangle edges respectively. We symmetrically define  $c_{I\Delta}$  and  $c_{\Delta\Delta}$ . If  $n_I(s)$  and  $n_\Delta(s)$  are the number of nodes of rank  $s$  which were infected along independent and triangle edges respectively, then

$$\begin{pmatrix} n_I(s+1) \\ n_\Delta(s+1) \end{pmatrix} = \begin{pmatrix} c_{II} & c_{I\Delta} \\ c_{\Delta I} & c_{\Delta\Delta} \end{pmatrix} \begin{pmatrix} n_I(s) \\ n_\Delta(s) \end{pmatrix},$$

where  $c_{II} = \frac{T\langle K_I^2 - K_I \rangle}{\langle K_I \rangle}$ ,  $c_{\Delta I} = \frac{2T(1+T-T^2)\langle K_I K_\Delta \rangle}{\langle K_I \rangle}$ ,  $c_{I\Delta} = \frac{T\langle K_I K_\Delta \rangle}{\langle K_\Delta \rangle}$ , and  $c_{\Delta\Delta} = \frac{2T(1+T-T^2)\langle K_\Delta^2 - K_\Delta \rangle}{\langle K_\Delta \rangle}$ . We give a sample calculation for  $c_{\Delta I}$ : With probability  $k_I p(k_I, k_\Delta) / \langle K_I \rangle$  an infection along an independent edge reaches a node with degrees  $k_I$  and  $k_\Delta$ . The expected number of infections along a triangle edge is  $2T(1+T-T^2)k_\Delta$ . Thus a random node infected along an independent edge creates  $2T(1+T-T^2)\langle K_I K_\Delta \rangle / \langle K_I \rangle$  infections along triangle edges.

The dominant eigenvalue of this matrix is  $\mathcal{R}_0$ . We generally want to determine  $T$  such that  $\mathcal{R}_0 < 1$ . Substituting  $\mathcal{R}_0 = 1$  into the characteristic equation

$$\begin{aligned} & \left( T \frac{\langle K_I^2 - K_I \rangle}{\langle K_I \rangle} - \mathcal{R}_0 \right) \left( 2T(1+T-T^2) \frac{\langle K_\Delta^2 - K_\Delta \rangle}{\langle K_\Delta \rangle} - \mathcal{R}_0 \right) \\ &= 2T^2(1+T-T^2) \frac{\langle K_I K_\Delta \rangle^2}{\langle K_I \rangle \langle K_\Delta \rangle} \end{aligned} \quad (1)$$

gives the critical transmissibility  $T = T_c$  below which epidemics have zero probability. At the threshold each factor on the left hand side is at most zero.

The original network has a giant component if  $\mathcal{R}_0 > 1$  when  $T = 1$ . Setting  $\mathcal{R}_0 = 1 + \mu$  and  $T = 1$  in (1), we let  $\chi(\mu)$  and  $\psi$  be the left and right hand side respectively. The concave parabola  $\chi$  has a negative minimum at  $\hat{\mu} = \langle K_I^2 - K_I \rangle / 2 \langle K_I \rangle + \langle K_\Delta^2 - K_\Delta \rangle / \langle K_\Delta \rangle - 1$ . Clearly  $\chi(\mu) > \psi$  as  $\mu \rightarrow \infty$ . If  $\hat{\mu} > 0$ , then the intermediate value theorem guarantees a positive root of

$\chi(\mu) = \psi$  (greater than  $\hat{\mu}$ ). Similarly, even if the first condition fails,  $\chi(0) < \psi$  guarantees a positive root. If both conditions fail there is no positive root. A giant component exists if

$$\frac{\langle K_I^2 - K_I \rangle}{2 \langle K_I \rangle} + \frac{\langle K_\Delta^2 - K_\Delta \rangle}{\langle K_\Delta \rangle} > 1$$

and/or

$$\left( \frac{\langle K_I^2 - K_I \rangle}{\langle K_I \rangle} - 1 \right) \left( 2 \frac{\langle K_\Delta^2 - K_\Delta \rangle}{\langle K_\Delta \rangle} - 1 \right) < 2 \frac{\langle K_I K_\Delta \rangle^2}{\langle K_I \rangle \langle K_\Delta \rangle}$$

If the first condition applies but not the second, then the network has enough independent and triangle edges that a giant component exists solely within the independent edges and another exists solely within the triangle edges. In all other cases the second condition applies.

### B. Unclustered, Segregated Network

We define  $n_b(s)$  and  $n_r(s)$  in the same manner, except that triangles need not be considered. Then

$$\begin{pmatrix} n_b(s+1) \\ n_r(s+1) \end{pmatrix} = \begin{pmatrix} c_{bb} & c_{br} \\ c_{rb} & c_{rr} \end{pmatrix} \begin{pmatrix} n_b \\ n_r \end{pmatrix}$$

where  $c_{bb} = \frac{T \langle K_b^2 - K_b \rangle}{\langle K_b \rangle}$ ,  $c_{br} = \frac{T \langle K_b K_r \rangle}{\langle K_r \rangle}$ ,  $c_{rb} = \frac{T \langle K_r K_b \rangle}{\langle K_b \rangle}$ , and  $c_{rr} = \frac{T \langle K_r^2 - K_r \rangle}{\langle K_r \rangle}$ . Substituting  $\mathcal{R}_0 = 1$  into the characteristic equation

$$\left( T \frac{\langle K_b^2 - K_b \rangle}{\langle K_b \rangle} - \mathcal{R}_0 \right) \left( T \frac{\langle K_r^2 - K_r \rangle}{\langle K_r \rangle} - \mathcal{R}_0 \right) = T^2 \frac{\langle K_r K_b \rangle^2}{\langle K_r \rangle \langle K_b \rangle} \quad (2)$$

finds the epidemic threshold. We divide (1) by  $1 + T - T^2$  and substitute  $\langle K_r^2 - K_r \rangle / \langle K_r \rangle = 2 \left( \langle K_\Delta^2 - K_\Delta \rangle / \langle K_\Delta \rangle \right) + 1$  and  $\langle K_r K_b \rangle^2 / \langle K_r \rangle \langle K_b \rangle = 2 \langle K_\Delta K_I \rangle / \langle K_\Delta \rangle \langle K_I \rangle$  into (2). Comparing the terms in the resulting equations shows that the threshold  $T$  in the unclustered network is at most the threshold in the corresponding clustered network.

### III. CALCULATING GIANT COMPONENT SIZE

To calculate the fraction of nodes in the giant component, it suffices to calculate the probability that a random node is *not* part of the giant component. These calculations have been done for CM networks by [11, 13, 16].

#### A. Clustered Network

We follow the approach of [13, 14]. A related approach is given by [16].

We let  $f$  be the probability a random node  $u$  is not part of the giant component. We have

$$f = \sum_{k_I, k_\Delta} p(k_I, k_\Delta) g_I^{k_I} g_\Delta^{k_\Delta},$$

where  $g_I$  and  $g_\Delta$  are the probabilities that an independent edge or a triangle respectively does not connect to the giant component. To find  $g_I$ , we note that there are two ways an edge can fail to connect  $u$  to the giant component: It may be deleted in the percolation process with probability  $1 - T$ , or it may be kept, but  $v$ , the node reached, is not part of the giant component. We have

$$g_I = 1 - T + T h_I$$

where  $h_I$  is the probability that a node  $v$  reached along an independent edge is not part of the giant component. To calculate  $h_I$  we note that  $v$  is selected proportional to  $k_I$ , but only has  $k_I - 1$  susceptible neighbors along independent edges. We get

$$h_I = \frac{1}{\langle K_I \rangle} \sum_{k_I, k_\Delta} k_I p(k_I, k_\Delta) g_I^{k_I - 1} g_\Delta^{k_\Delta}.$$

For  $g_\Delta$  we get

$$g_\Delta = [1 - T + T h_\Delta]^2 - 2T^2(1 - T)h_\Delta(1 - h_\Delta),$$

where  $h_\Delta$  is the probability a node reached along a triangle edge does not connect to the giant component through any edge not in the triangle. We find

$$h_\Delta = \frac{1}{\langle K_\Delta \rangle} \sum_{k_I, k_\Delta} k_\Delta p(k_I, k_\Delta) g_I^{k_I} g_\Delta^{k_\Delta - 1}.$$

The resulting system of equations for  $g_I$ ,  $g_\Delta$ ,  $h_I$ , and  $h_\Delta$  can be solved iteratively, and the result gives  $f$ .

#### B. Unclustered, Segregated Network

To find  $f_u$ , the probability a random node in the unclustered network is not part of the giant component, we proceed similarly. We find

$$\begin{aligned} f_u &= \sum_{k_r, k_b} p_u(k_b, k_r) g_b^{k_b} g_r^{k_r} \\ g_b &= 1 - T + T h_b \\ g_r &= 1 - T + T h_r \\ h_b &= \frac{1}{\langle K_b \rangle} \sum_{k_r, k_b} k_b p_u(k_b, k_r) g_b^{k_b - 1} g_r^{k_r} \\ h_r &= \frac{1}{\langle K_r \rangle} \sum_{k_r, k_b} k_r p_u(k_b, k_r) g_b^{k_b} g_r^{k_r - 1}. \end{aligned}$$

By iterating beginning with  $h_b$  and  $h_r$  both zero we find  $g_b$  and  $g_r$ , from which  $f_u$  can be calculated. A similar approach will find  $f$  for corresponding clustered networks. At each step of the iteration,  $g_r^2 \leq g_\Delta$  and  $g_b \leq g_I$ ,

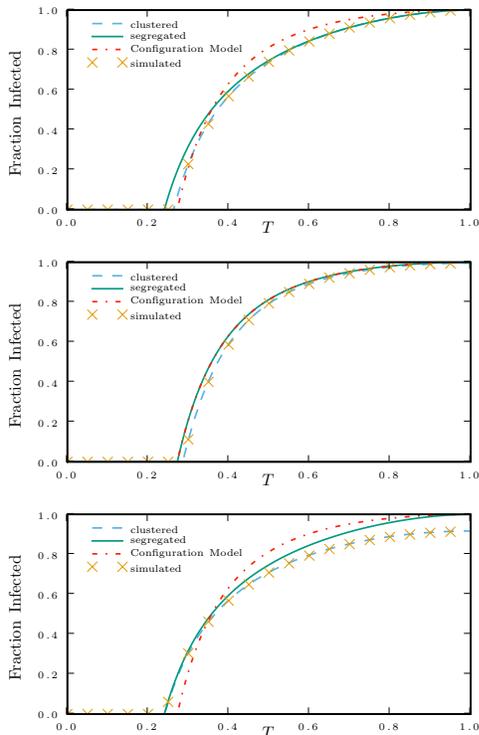


FIG. 1: A comparison of different network configurations. Assortative mixing reduces the epidemic threshold. Clustering reduces epidemic size.

from which we conclude  $f_u \leq f$ . Consequently the size of the giant component is smaller in clustered networks than in unclustered networks of the same degree distribution and degree correlations.

#### IV. RESULTS

In figure 1 we consider outbreak spread on three networks, all of which have the same degree distribution. We compare simulated epidemic sizes with predictions from the clustered equations, the unclustered, segregated equations, and the equations derived previously for configuration model networks [11, 13, 16].

The nodes are equally distributed between degrees 2, 4, and 6. In each network the clustering is distributed differently. In the first,  $p(0,3) = 1/3$ ,  $p(2,1) = 1/3$ , and  $p(2,0) = 1/3$ . That is those nodes with degree 6 are only in triangles, nodes of degree 4 have half of their edges in

triangles and independent edges, and nodes of degree 2 have just independent edges. High degree nodes tend to be clustered and contact other high degree nodes. The tendency to contact other high degree nodes reduces the epidemic threshold, but the clustering raises the threshold.

In the second network, we take  $p(2,0) = 1/6$ ,  $p(0,1) = 1/6$ ,  $p(2,1) = 1/3$ ,  $p(4,1) = 1/6$ , and  $p(0,3) = 1/6$ . This yields identical distribution of neighbor degrees for nodes reached by either a triangle or an independent edge. The unclustered, segregated equations yield the same result as the configuration model equations. The clustered calculations have smaller epidemics.

The third network is an inversion of the first. Nodes with high degree have independent edges while nodes with low degree are clustered. We take  $p(6,0) = 1/3$ ,  $p(2,1) = 1/3$ , and  $p(0,1) = 1/3$ . Again the assortativity reduces the epidemic threshold while clustering reduces the epidemic size. In this particular case, it is the preference for high degree nodes (which are unclustered) to contact one another that leads to the reduction in epidemic threshold, and so it is clear that the effect is due to assortative mixing, not clustering.

#### V. DISCUSSION

We have introduced a new model of clustered networks on which we study percolation and epidemics. This model allows us to make a number of analytic prediction because the edges of the network can be partitioned into sets which are independent of one another (independent edges or triangles).

We have shown that these networks can have a lower epidemic threshold than Configuration Model networks with the same degree distribution. However, this is not a consequence of clustering, but rather a consequence of assortative mixing. The clustering of the network can be proven to raise the epidemic threshold and reduce the epidemic size from networks with the same degree correlations, but without clustering.

#### Acknowledgments

This work was supported by the RAPIDD program of the Science & Technology Directorate, Department of Homeland Security and the Fogarty International Center, National Institutes of Health.

- 
- [1] Shweta Bansal. *Ecology of Infectious Diseases with Contact Networks and Percolation Theory*. PhD thesis, University of Texas at Austin, 2008.  
 [2] N.G. Becker, K. Glass, Z. Li, and G.K. Aldis. Controlling

emerging infectious diseases like SARS. *Mathematical biosciences*, 193(2):205–221, 2005.

- [3] B. Bollobás. A probabilistic proof of an asymptotic formula for the number of labelled random graphs. *European*

- Journal of Combinatorics*, 1:311–316, 1980.
- [4] T. Britton, M. Deijfen, A.N. Lageras, and M. Lindholm. Epidemics on random graphs with tunable clustering. *Journal of Applied Probability*, 45:743–756, 2008.
- [5] John L. Cardy and Peter Grassberger. Epidemic models and percolation. *Journal of Physics A: Mathematics and General*, 18(6):L267–L271, 1985.
- [6] O. Diekmann, J. A. P. Heesterbeek, and J. A. J. Metz. On the definition and the computation of the basic reproduction ratio  $\mathcal{R}_0$  in models for infectious diseases in heterogeneous populations. *Journal of Mathematical Biology*, 28:365–382, 1990.
- [7] K. T. D. Eames. Modelling disease spread through random and regular contacts in clustered populations. *Theoretical Population Biology*, 73:104–111, 2008.
- [8] Peter Grassberger. On the critical behavior of the general epidemic process and dynamical percolation. *Mathematical Biosciences*, 63:157–172, 1983.
- [9] M. J. Keeling. The effects of local spatial structure on epidemiological invasions. *Proceedings of the Royal Society B: Biological Sciences*, 266(1421):859–867, 1999.
- [10] Eben Kenah and James M. Robins. Network-based analysis of stochastic SIR epidemic models with random and proportionate mixing. *Journal of Theoretical Biology*, 249(4):706–722, 2007.
- [11] Eben Kenah and James M. Robins. Second look at the spread of epidemics on networks. *Physical Review E*, 76(3):36113, 2007.
- [12] D. Ludwig. Final size distributions for epidemics. *Mathematical Biosciences*, 23:33–46, 1975.
- [13] Joel C. Miller. Epidemic size and probability in populations with heterogeneous infectivity and susceptibility. *Physical Review E*, 76(1):010101, 2007.
- [14] Joel C. Miller. Spread of infectious disease through clustered populations. *Journal of the Royal Society, Interface*, pages ??–??, 2009.
- [15] M. Molloy and Bruce Reed. A critical point for random graphs with a given degree sequence. *Random structures & algorithms*, 6(2):161–179, 1995.
- [16] Mark E. J. Newman. Spread of epidemic disease on networks. *Physical Review E*, 66(1):16128, 2002.
- [17] Mark E. J. Newman. Properties of highly clustered networks. *Physical Review E*, 68(2):026121, 2003.
- [18] Mark E. J. Newman. The structure and function of complex networks. *SIAM Review*, 45:167–256, 2003.
- [19] MEJ Newman. Random graphs with clustering. *Arxiv preprint arXiv:0903.4009*, 2009.
- [20] M. Ángeles Serrano and Marián Boguñá. Clustering in complex networks. II. Percolation properties. *Physical Review E*, 74(5):056115, 2006.
- [21] M. Ángeles Serrano and Marián Boguñá. Percolation and epidemic thresholds in clustered networks. *Physical Review Letters*, 97(8):088701, 2006.
- [22] Care must be taken that no dependence between edges arises. Such a dependence can arise from, for example, heterogeneity in duration of infection [11, 13].
- [23] This algorithm was simultaneously developed by [19]
- [24] In essence we have created a generalization of an edge which corresponds to a triangle. We could create other more general structures in much the same way
- [25] Without our simplification, we would need to further subdivide those infected along triangle edges into those whose other neighbor is still susceptible from those whose other neighbor is also infected.